

CS 726: Homework #4

Posted: 03/08/2020, due: 03/27/2020 by 5pm on Canvas

Please typeset or write your solutions neatly! If we cannot read it, we cannot grade it.

Note: You can use the results we have proved in class – no need to prove them again.

Q 1. Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be an L -smooth function and let $\mathcal{X} \subseteq \mathbb{R}^n$ be a closed, convex, and nonempty set. Recall the definition of the gradient mapping: $G_\eta(\mathbf{x}) = \eta(\mathbf{x} - P_{\mathcal{X}}(\mathbf{x} - \frac{1}{L}\nabla f(\mathbf{x})))$, where $P_{\mathcal{X}}(\cdot)$ denotes the Euclidean projection onto \mathcal{X} . Prove that $G_\eta(\cdot)$ is $(2\eta + L)$ -Lipschitz continuous. [20pts]

Q 2. Consider a constrained minimization problem $\min_{\mathbf{u} \in \mathcal{X}} f(\mathbf{x})$, where $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is L -smooth and $\mathcal{X} \subseteq \mathbb{R}^n$ is a hyper-rectangle: $\mathcal{X} = \{\mathbf{x} \in \mathbb{R}^n : x_i \in [a_i, b_i], \forall i \in \{1, 2, \dots, n\}\}$, where $\mathbf{a}, \mathbf{b} \in \mathbb{R}^n$ are vectors that satisfy $\forall i \in \{1, \dots, n\} : a_i < b_i$.

- (i) What is $P_{\mathcal{X}}(\mathbf{x})$? Write it in closed form. [10pts]
- (ii) Define $\Delta_i(\mathbf{x}) := \nabla_i f(\mathbf{x}) \mathbf{e}_i$, where \mathbf{e}_i is the i^{th} standard basis vector (except for its i^{th} coordinate, all coordinates are equal to zero; the i^{th} coordinate equals one). Define:

$$T(\mathbf{x}, i) := \underset{\mathbf{u} \in \mathcal{X}}{\operatorname{argmin}} \left\{ \langle \Delta_i(\mathbf{x}), \mathbf{u} - \mathbf{x} \rangle + \frac{L}{2} \|\mathbf{u} - \mathbf{x}\|_2^2 \right\}.$$

Express the gradient mapping $G_L(\mathbf{x})$ as a function of L , \mathbf{x} , and $T(\mathbf{x}, i)$, $i \in \{1, \dots, n\}$. [10pts]

- (iii) Consider the following method that starts from some initial point $\mathbf{x}_0 \in \mathcal{X}$ and updates its iterates \mathbf{x}_k for $k \geq 0$ as:

$$\begin{aligned} i_k^* &= \underset{1 \leq i \leq n}{\operatorname{argmax}} |(G_L(\mathbf{x}_k))_i| \\ \mathbf{x}_{k+1} &= T(\mathbf{x}_k, i_k^*). \end{aligned}$$

Prove the following sufficient descent property for this algorithm:

$$(\exists \alpha > 0)(\forall k \geq 0) : f(\mathbf{x}_{k+1}) \leq f(\mathbf{x}_k) - \frac{\alpha}{2} \|G_L(\mathbf{x}_k)\|_2^2.$$

What is the largest α for which this property holds? What can you say about convergence of this method if f is bounded below by some \bar{f} ? [10pts]

Q 3. Consider a constrained minimization problem $\min_{\mathbf{u} \in \mathcal{X}} f(\mathbf{x})$, where $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is convex and M -Lipschitz, and $\mathcal{X} \subseteq \mathbb{R}^n$ is closed, convex, and nonempty. Assume that there exists $\mathbf{x}^* \in \mathcal{X}$ that minimizes f . Recall from the class that we use \mathbf{g}_x to denote an arbitrary subgradient of f at $\mathbf{x} \in \mathcal{X}$. Distances within the set \mathcal{X} are measured using some norm $\|\cdot\|$ (e.g., an ℓ_p norm for $p \geq 1$). Norm that is dual to $\|\cdot\|$ is denoted by $\|\cdot\|_*$ and is defined by: $\|\mathbf{z}\|_* = \max_{\mathbf{u}: \|\mathbf{u}\| \leq 1} \langle \mathbf{z}, \mathbf{u} \rangle$. By the definition of a dual norm, you can immediately deduce that:

$$(\forall \mathbf{u}, \mathbf{z}) : \langle \mathbf{u}, \mathbf{z} \rangle \leq \|\mathbf{u}\| \|\mathbf{z}\|_*.$$

Note also that since f is M -Lipschitz w.r.t. the norm $\|\cdot\|$, we have, $\forall \mathbf{x}, \mathbf{y}$ and all $\mathbf{g}_x \in \partial f(\mathbf{x})$:

$$|f(\mathbf{x}) - f(\mathbf{y})| \leq M \|\mathbf{x} - \mathbf{y}\| \quad \text{and } \|\mathbf{g}_x\|_* \leq M.$$

In this part, you will analyze the convergence of the Mirror Descent algorithm, defined by:

$$\mathbf{x}_{k+1} = \underset{\mathbf{u} \in \mathcal{X}}{\operatorname{argmin}} \{ \langle \mathbf{g}_{\mathbf{x}_k}, \mathbf{u} - \mathbf{x}_k \rangle + D_\psi(\mathbf{u}, \mathbf{x}_k) \}, \tag{MD}$$

where, as in previous assignments, $D_\psi(\mathbf{x}, \mathbf{y})$ denotes the Bregman divergence between \mathbf{x} and \mathbf{y} w.r.t. a continuously-differentiable function ψ .

Assume that ψ is 1-strongly convex w.r.t. $\|\cdot\|$, namely:

$$(\forall \mathbf{x}, \mathbf{y}) : \quad \psi(\mathbf{y}) \geq \psi(\mathbf{x}) + \langle \nabla \psi(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + \frac{1}{2} \|\mathbf{y} - \mathbf{x}\|^2.$$

- (i) Using the first-order optimality in the definition of \mathbf{x}_{k+1} and the three-point identity of Bregman divergences we proved in HW #2, show that:

$$(\forall \mathbf{u} \in \mathcal{X}) : \quad a_k \langle \nabla f(\mathbf{x}_k), \mathbf{x}_{k+1} - \mathbf{u} \rangle \leq D_\psi(\mathbf{x}^*, \mathbf{x}_k) - D_\psi(\mathbf{x}^*, \mathbf{x}_{k+1}) - D_\psi(\mathbf{x}_{k+1}, \mathbf{x}_k). \quad [15\text{pts}]$$

- (ii) Use Part (i) to prove that:

$$a_k \langle \nabla f(\mathbf{x}_k), \mathbf{x}_k - \mathbf{x}^* \rangle \leq \frac{a_k^2 M^2}{2} + D_\psi(\mathbf{x}^*, \mathbf{x}_k) - D_\psi(\mathbf{x}^*, \mathbf{x}_{k+1}). \quad [15\text{pts}]$$

- (iii) Use Part (ii) to prove that, $\forall k \geq 0$:

$$f(\mathbf{x}_k^{\text{out}}) - f(\mathbf{x}^*) \leq \frac{D_\psi(\mathbf{x}^*, \mathbf{x}_0) + \frac{M^2}{2} \sum_{i=0}^k a_i^2}{A_k},$$

$$\text{where } \mathbf{x}_k^{\text{out}} = \frac{\sum_{i=0}^k a_i \mathbf{x}_i}{A_k}. \quad [10\text{pts}]$$

- (iv) Discuss how you would choose $\{a_i\}_{i=0}^k$ for (MD) to converge as fast as possible, and provide the convergence bound (a bound on $f(\mathbf{x}_k^{\text{out}}) - f(\mathbf{x}^*)$) for your choice(s) of $\{a_i\}_{i=0}^k$. [10pts]

Q 1. Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be an L -smooth function and let $\mathcal{X} \subseteq \mathbb{R}^n$ be a closed, convex, and nonempty set. Recall the definition of the gradient mapping: $G_\eta(\mathbf{x}) = \eta(\mathbf{x} - P_{\mathcal{X}}(\mathbf{x} - \frac{1}{\eta} \nabla f(\mathbf{x})))$, where $P_{\mathcal{X}}(\cdot)$ denotes the Euclidean projection onto \mathcal{X} . Prove that $G_\eta(\cdot)$ is $(2\eta + L)$ -Lipschitz continuous. [20pts]

$\forall \mathbf{x}, \mathbf{y} \in \mathcal{X}$.

$$G_\eta(\mathbf{x}) - G_\eta(\mathbf{y}) = \eta \left[\mathbf{x} - \mathbf{y} - \left(P_{\mathcal{X}}\left(\mathbf{x} - \frac{1}{\eta} \nabla f(\mathbf{x})\right) - P_{\mathcal{X}}\left(\mathbf{y} - \frac{1}{\eta} \nabla f(\mathbf{y})\right) \right) \right]$$

$$\begin{aligned} |G_\eta(\mathbf{x}) - G_\eta(\mathbf{y})| &\leq \eta \|\mathbf{x} - \mathbf{y}\| + \eta \left\| P_{\mathcal{X}}\left(\mathbf{x} - \frac{1}{\eta} \nabla f(\mathbf{x})\right) - P_{\mathcal{X}}\left(\mathbf{y} - \frac{1}{\eta} \nabla f(\mathbf{y})\right) \right\| \\ &\leq \eta \|\mathbf{x} - \mathbf{y}\| + \eta \left\| \mathbf{x} - \mathbf{y} + \frac{1}{\eta} (\nabla f(\mathbf{y}) - \nabla f(\mathbf{x})) \right\| \\ &\leq 2\eta \|\mathbf{x} - \mathbf{y}\| + \|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\| \\ &\leq (2\eta + L) \|\mathbf{x} - \mathbf{y}\| \end{aligned}$$

Q 2. Consider a constrained minimization problem $\min_{\mathbf{u} \in \mathcal{X}} f(\mathbf{x})$, where $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is L -smooth and $\mathcal{X} \subseteq \mathbb{R}^n$ is a hyper-rectangle: $\mathcal{X} = \{\mathbf{x} \in \mathbb{R}^n : x_i \in [a_i, b_i], \forall i \in \{1, 2, \dots, n\}\}$, where $\mathbf{a}, \mathbf{b} \in \mathbb{R}^n$ are vectors that satisfy $\forall i \in \{1, \dots, n\} : a_i < b_i$.

(i) What is $P_{\mathcal{X}}(\mathbf{x})$? Write it in closed form. $\min \left\{ \max \{a_i, x_i\}, b_i \right\}$ [10pts]

(ii) Define $\Delta_i(\mathbf{x}) := \nabla_i f(\mathbf{x}) \mathbf{e}_i$, where \mathbf{e}_i is the i^{th} standard basis vector (except for its i^{th} coordinate, all coordinates are equal to zero; the i^{th} coordinate equals one). Define:

$$T(\mathbf{x}, i) := \underset{\mathbf{u} \in \mathcal{X}}{\operatorname{argmin}} \left\{ \langle \Delta_i(\mathbf{x}), \mathbf{u} - \mathbf{x} \rangle + \frac{L}{2} \|\mathbf{u} - \mathbf{x}\|_2^2 \right\}.$$

Express the gradient mapping $G_L(\mathbf{x})$ as a function of L , \mathbf{x} , and $T(\mathbf{x}, i)$, $i \in \{1, \dots, n\}$. [10pts]

(iii) Consider the following method that starts from some initial point $\mathbf{x}_0 \in \mathcal{X}$ and updates its iterates \mathbf{x}_k for $k \geq 0$ as:

$$\begin{aligned} i_k^* &= \underset{1 \leq i \leq n}{\operatorname{argmax}} |(G_L(\mathbf{x}_k))_i| \\ \mathbf{x}_{k+1} &= T(\mathbf{x}_k, i_k^*). \end{aligned}$$

Prove the following sufficient descent property for this algorithm:

$$(\exists \alpha > 0)(\forall k \geq 0) : f(\mathbf{x}_{k+1}) \leq f(\mathbf{x}_k) - \frac{\alpha}{2} \|G_L(\mathbf{x}_k)\|_2^2.$$

What is the largest α for which this property holds? What can you say about convergence of this method if f is bounded below by some \bar{f} ? [10pts]

$$(i) P_{\mathcal{X}}(\mathbf{x}) = \min \left\{ \max \{a_i, x_i\}, b_i \right\}$$

$$(ii) \text{ Let } \bar{\mathbf{x}} = \underset{\mathbf{u} \in \mathcal{X}}{\operatorname{argmin}} \left\{ \langle \nabla f(\mathbf{x}), \mathbf{u} - \mathbf{x} \rangle + \frac{L}{2} \|\mathbf{u} - \mathbf{x}\|_2^2 \right\} = P_{\mathcal{X}}\left(\mathbf{x} - \frac{1}{L} \nabla f(\mathbf{x})\right) = \mathbf{x} - \frac{1}{L} G_L(\mathbf{x})$$

$$G_L(\mathbf{x}) = L \left(\mathbf{x} - P_{\mathcal{X}}\left(\mathbf{x} - \frac{1}{L} \nabla f(\mathbf{x})\right) \right)$$

We have:

$$T(x, i) = \underset{u}{\operatorname{argmin}} \left\{ \langle \Delta_i(x), u - x \rangle + \frac{L}{2} \|u - x\|_v^2 \right\} = P_X \left(x - \frac{\Delta_i(x)}{L} \right) = x - \frac{1}{L} G_L(x) \cdot e_i$$

$$\text{define } G_L^i(x) = G_L(x) \cdot e_i$$

$$= L(x - T(x, i))$$

$$G_L(x) = Lx - L \cdot \begin{bmatrix} T(x, 1) \\ \vdots \\ T(x, n) \end{bmatrix}$$

$$(iii) \quad x_{k+1} = x_k - \frac{1}{L} G_L^{i_k}(x_k) \quad \text{where } i_k = \underset{0 \leq i \leq n}{\operatorname{argmax}} |[G_L(x)]_i|$$

$$\begin{aligned} f(x_{k+1}) &\leq f(x_k) + \langle \nabla f(x_k), x_{k+1} - x_k \rangle + \frac{L}{2} \|x_{k+1} - x_k\|_v^2 \\ &\leq f(x_k) - \frac{1}{L} \langle \nabla f(x_k), G_L^{i_k}(x_k) \rangle + \frac{1}{2L} \|G_L^{i_k}(x_k)\|_v^2 \\ &= f(x_k) - \frac{1}{L} \langle \Delta_{i_k}(x_k), G_L^{i_k}(x_k) \rangle + \frac{1}{2L} \|G_L^{i_k}(x_k)\|_v^2 \\ &= f(x_k) - \frac{1}{L} \langle \Delta_{i_k}(x_k) - G_L^{i_k}(x_k), G_L^{i_k}(x_k) \rangle - \frac{1}{2L} \|G_L^{i_k}(x_k)\|_v^2 \end{aligned}$$

$$\begin{aligned} \text{Since: } &\langle \Delta_{i_k}(x_k) - G_L^{i_k}(x_k), G_L^{i_k}(x_k) \rangle \\ &= L^2 \langle x_k - \left(x_k - \frac{1}{L} \Delta_{i_k}(x_k) \right) + T(x_k, i_k), x_k - T(x_k, i_k) \rangle = 1 \end{aligned}$$

$$\text{Consider: } T(x, i_k) = P_X(x - \frac{1}{L} \Delta_{i_k}(x))$$

According to the property of Euclidean projection:

$$\langle x - P_X(x), u - P_X(x) \rangle \leq 0 \quad \forall u \in X$$

$$\text{Then: } 1 = L^2 \langle P_X(x - \frac{1}{L} \Delta_{i_k}(x)) - (x - \frac{1}{L} \Delta_{i_k}(x)), x - P_X(x - \frac{1}{L} \Delta_{i_k}(x)) \rangle \geq 0$$

$$\text{We have: } f(x_{k+1}) \leq f(x_k) - \frac{1}{2L} \|G_L(x_k) \cdot e_{i_k}\|^2$$

$$\leq f(x_k) - \frac{1}{2nL} \|G_L(x_k)\|^2$$

$$\lambda = \frac{1}{nL}$$

If $f(x) \geq \bar{f} \quad \forall x$

$$f(x_{k+1}) - f(x_0) \leq -\frac{1}{2nL} \sum_{i=0}^k \|G_L(x_i)\|^2$$

$$\min_{0 \leq i \leq k} \|G_L(x_i)\|^2 \leq \frac{2nL}{k+1} (f(x_0) - f(x_{k+1})) \leq \frac{2nL}{k+1} (f(x_0) - \bar{f})$$

Q 3. Consider a constrained minimization problem $\min_{\mathbf{u} \in \mathcal{X}} f(\mathbf{x})$, where $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is convex and M -Lipschitz, and $\mathcal{X} \subseteq \mathbb{R}^n$ is closed, convex, and nonempty. Assume that there exists $\mathbf{x}^* \in \mathcal{X}$ that minimizes f . Recall from the class that we use \mathbf{g}_x to denote an arbitrary subgradient of f at $\mathbf{x} \in \mathcal{X}$. Distances within the set \mathcal{X} are measured using some norm $\|\cdot\|$ (e.g., an ℓ_p norm for $p \geq 1$). Norm that is dual to $\|\cdot\|$ is denoted by $\|\cdot\|_*$ and is defined by: $\|\mathbf{z}\|_* = \max_{\mathbf{u}: \|\mathbf{u}\| \leq 1} \langle \mathbf{z}, \mathbf{u} \rangle$. By the definition of a dual norm, you can immediately deduce that:

$$(\forall \mathbf{u}, \mathbf{z}) : \langle \mathbf{u}, \mathbf{z} \rangle \leq \|\mathbf{u}\| \|\mathbf{z}\|_*$$

Note also that since f is M -Lipschitz w.r.t. the norm $\|\cdot\|$, we have, $\forall \mathbf{x}, \mathbf{y}$ and all $\mathbf{g}_x \in \partial f(\mathbf{x})$:

$$|f(\mathbf{x}) - f(\mathbf{y})| \leq M \|\mathbf{x} - \mathbf{y}\| \quad \text{and } \|\mathbf{g}_x\|_* \leq M.$$

In this part, you will analyze the convergence of the Mirror Descent algorithm, defined by:

$$\begin{aligned} \mathbf{x}_{k+1} &= \underset{\mathbf{u} \in \mathcal{X}}{\operatorname{argmin}} \{ \langle \mathbf{g}_{x_k}, \mathbf{u} - \mathbf{x}_k \rangle + D_\psi(\mathbf{u}, \mathbf{x}_k) \}, \\ \mathbf{x}_{k+1} &= \underset{\mathbf{u} \in \mathcal{X}}{\operatorname{argmin}} \left\{ a_k \langle \mathbf{g}_{x_k}, \mathbf{u} - \mathbf{x}_k \rangle + D_\psi(\mathbf{u}, \mathbf{x}_k) \right\} \end{aligned} \quad \begin{aligned} &\leftarrow a_k \langle \mathbf{g}_{x_k} - \nabla f(x_{k+1}), \mathbf{u} - \mathbf{x}_{k+1} \rangle \leq 0 \\ &\leftarrow a_k \langle \mathbf{g}_{x_k} - \nabla f(x_{k+1}), \mathbf{x}_{k+1} - \mathbf{x}_k \rangle \leq 0 \end{aligned} \quad \begin{aligned} &\forall u \\ &\forall u \end{aligned}$$

where, as in previous assignments, $D_\psi(\mathbf{x}, \mathbf{y})$ denotes the Bregman divergence between \mathbf{x} and \mathbf{y} w.r.t. a continuously-differentiable function ψ .

Assume that ψ is 1-strongly convex w.r.t. $\|\cdot\|$, namely:

$$(\forall \mathbf{x}, \mathbf{y}) : \psi(\mathbf{y}) \geq \psi(\mathbf{x}) + \langle \nabla \psi(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + \frac{1}{2} \|\mathbf{y} - \mathbf{x}\|^2.$$

- (i) Using the first-order optimality in the definition of \mathbf{x}_{k+1} and the three-point identity of Bregman divergences we proved in HW #2, show that:

$$\langle \nabla \psi(x_{k+1}) - \nabla \psi(x_k), \mathbf{x} - \mathbf{x}_{k+1} \rangle >$$

$$(\forall \mathbf{u} \in \mathcal{X}) : a_k \langle \nabla f(\mathbf{x}_k), \mathbf{x}_{k+1} - \mathbf{u} \rangle \leq D_\psi(\mathbf{x}^*, \mathbf{x}_k) - D_\psi(\mathbf{x}^*, \mathbf{x}_{k+1}) - D_\psi(\mathbf{x}_{k+1}, \mathbf{x}_k). \quad [15\text{pts}]$$

- (ii) Use Part (i) to prove that:

$$a_k \langle \nabla f(\mathbf{x}_k), \mathbf{x}_k - \mathbf{x}^* \rangle \leq \frac{a_k^2 M^2}{2} + D_\psi(\mathbf{x}^*, \mathbf{x}_k) - D_\psi(\mathbf{x}^*, \mathbf{x}_{k+1}). \quad [15\text{pts}]$$

- (iii) Use Part (ii) to prove that, $\forall k \geq 0$:

$$f(\mathbf{x}_k^{\text{out}}) - f(\mathbf{x}^*) \leq \frac{D_\psi(\mathbf{x}^*, \mathbf{x}_0) + \frac{M^2}{2} \sum_{i=0}^k a_i^2}{A_k},$$

$$\text{where } \mathbf{x}_k^{\text{out}} = \frac{\sum_{i=0}^k a_i \mathbf{x}_i}{A_k}. \quad [10\text{pts}]$$

- (iv) Discuss how you would choose $\{a_i\}_{i=0}^k$ for (MD) to converge as fast as possible, and provide the convergence bound (a bound on $f(\mathbf{x}_k^{\text{out}}) - f(\mathbf{x}^*)$) for your choice(s) of $\{a_i\}_{i=0}^k$. [10pts]

$$(i) \text{ Since } X_{k+1} = \underset{u \in \mathcal{X}}{\operatorname{arg\min}} \left\{ \alpha_k \langle g_{X_k}, u - X_k \rangle + D_\Psi(u, X_k) \right\}$$

$$D_\Psi(x, y) = \Psi(x) - \Psi(y) - \langle \nabla \Psi(y), x - y \rangle$$

$$\begin{aligned} \text{we have } & \langle -\alpha_k g_{X_k} - \nabla D_\Psi(u, X_k) \Big|_{u=X_{k+1}}, u - X_{k+1} \rangle \leq 0 \quad \forall u \in \mathcal{X} \\ & - \langle \alpha_k g_{X_k} + \nabla \Psi(X_{k+1}) - \nabla \Psi(X_k), X_{k+1} - u \rangle \geq 0 \\ & \alpha_k \langle g_{X_k}, X_{k+1} - u \rangle \leq \langle \nabla \Psi(X_{k+1}) - \nabla \Psi(X_k), u - X_{k+1} \rangle \quad \forall u \in \mathcal{X} \quad \text{①} \end{aligned}$$

According to three-point identity of Bregman-dissimilarity:

$$\langle \nabla \Psi(X_{k+1}) - \nabla \Psi(X_k), u - X_{k+1} \rangle = D_\Psi(u, X_k) - D_\Psi(X_{k+1}, X_k) - D_\Psi(u, X_{k+1})$$

where $u \in \mathcal{X}$ ②

$$\text{Since } f(y) \geq f(x) + \langle g_x, y - x \rangle. \text{ let } y = x + \lambda p \quad \text{if there } \|p\|=1$$

$$\langle g_x, p \rangle \leq \frac{f(x+\lambda p) - f(x)}{\lambda} = \langle \nabla f(x), p \rangle \quad \text{if } \nabla f(x) \text{ exists}$$

$$\langle g_x - \nabla f(x), p \rangle \leq 0 \quad \text{for all } p.$$

$$\text{then if } \nabla f(x) \text{ exists, } g_x = \nabla f(x) \quad \text{③}$$

Then combine ①. ②. ③.

$$\alpha_k \langle \nabla f(X_k), X_{k+1} - u \rangle \leq D_\Psi(X^*, X_k) - D_\Psi(X_{k+1}, X_k) - D_\Psi(X^*, X_{k+1}) \quad \forall u \in \mathcal{X}$$

(ii) According to (i). we have:

$$\alpha_k \langle \nabla f(X_k), X_{k+1} - X^* \rangle \leq D_\Psi(X^*, X_k) - D_\Psi(X_{k+1}, X_k) - D_\Psi(X^*, X_{k+1}) \quad \text{①}$$

According to $\Psi(x)$ is 1-strongly convex: $\forall y, x \in \mathcal{X}$

$$\begin{aligned} f(y) - f(x) - \langle \nabla f(x), y - x \rangle &\geq \frac{1}{2} \|y - x\|_v^2 \\ D_\Psi(y, x) &\geq \frac{1}{2} \|y - x\|_v^2 \quad \text{②} \end{aligned}$$

Then we have:

$$\begin{aligned} \alpha_k \langle \nabla f(X_k), X_k - X^* \rangle &\leq D_\Psi(X^*, X_k) - D_\Psi(X^*, X_{k+1}) - \frac{1}{2} \|X_{k+1} - X_k\|_v^2 \\ &\quad + \alpha_k \langle \nabla f(X_k), X_k - X_{k+1} \rangle \quad \text{③} \end{aligned}$$

$$\text{Claim: } -\frac{1}{2} \|x_{k+1} - x_k\|_v^2 + \alpha_k \langle \nabla f(x_k), x_k - x_{k+1} \rangle \leq \frac{\alpha_k^2 M^2}{2}$$

$$\begin{aligned} \text{pf: } \alpha_k \langle \nabla f(x_k), x_k - x_{k+1} \rangle &= \alpha_k \langle g_{x_k}, x_k - x_{k+1} \rangle \leq \alpha_k \|g_{x_k}\| \cdot \|x_k - x_{k+1}\| \\ &\leq \alpha_k M \cdot \|x_k - x_{k+1}\| \end{aligned}$$

$$\text{let } p = \|x_{k+1} - x_k\| \quad q = \alpha_k M$$

$$-\frac{1}{2} \|x_{k+1} - x_k\|_v^2 + \alpha_k \langle \nabla f(x_k), x_k - x_{k+1} \rangle - \frac{p^2}{2} + pq \leq \frac{q^2}{2} = \frac{\alpha_k^2 M^2}{2}$$

then complete the proof of claim.

Then complete the proof of iii). \square

(iii) Consider optimality gap:

$$U_k = \frac{1}{A_k} \sum_{i=0}^k \alpha_i f(x_i) \geq f\left(\frac{1}{A_k} \sum_{i=0}^k \alpha_i x_i\right) = f(x^{int})$$

$$\alpha_k U_k - \alpha_{k-1} U_{k-1} = \alpha_k f(x_k)$$

$$f(x^*) \geq \frac{1}{A_k} \sum_{i=0}^k \alpha_i \left(f(x_i) + \langle g_{x_i}, x^* - x_i \rangle \right) = L_k$$

$$\alpha_k L_k - \alpha_{k-1} L_{k-1} = \alpha_k f(x_k) - \alpha_k \langle g_{x_k}, x^* - x_k \rangle$$

$$\begin{aligned} \text{Then: } \alpha_k G_k - \alpha_{k-1} G_{k-1} &= -\alpha_k \langle g_{x_k}, x^* - x_k \rangle \\ &\leq \frac{\alpha_k^2 M^2}{2} + D_\psi(x^*, x_k) - D_\psi(x^*, x_{k+1}) \end{aligned}$$

$$\alpha_0 G_0 = \alpha_0 f(x_0) - \alpha_0 f(x_0) - \alpha_0 \langle g_{x_0}, x^* - x_0 \rangle \leq \frac{\alpha_0^2 M^2}{2} D_\psi(x^*, x_0) - D_\psi(x^*, x_1)$$

Thus:

$$\alpha_k G_k \leq \frac{M^2}{2} \sum_{i=0}^k \alpha_i^2 + D_\psi(x^*, x_0) - D_\psi(x^*, x_{k+1})$$

$$f(x^{int}) - f(x^*) \leq G_k \leq \underbrace{D_\psi(x^*, x_0) + \frac{\frac{M^2}{2} \sum_{i=0}^k \alpha_i^2}{A_k}}_{A_k}$$

(iv): Choose $\alpha_i = 2$

$$G_k \leq \frac{D\psi(x^*, x_0) + M^2 \alpha^2 (k+1)}{2\gamma(k+1)} = \frac{D\psi(x^*, x_0)}{2\gamma(k+1)} + \frac{2M^2}{z}$$

$$\text{Choose } 2 \text{ s.t. } \frac{D\psi(x^*, x_0)}{2\gamma(k+1)} = \frac{2M^2}{z} \Rightarrow 2 = \sqrt{\frac{D\psi(x^*, x_0)}{M^2(k+1)}}$$

$$\text{Then: } G_k \leq \frac{2M^2}{z} \\ = \sqrt{\frac{D\psi(x^*, x_0)}{k+1}} \cdot M \leq \zeta$$

$$f(x_k^{(m)}) - f(x^*) \leq \zeta \quad \text{in}$$

$$k = 0 \left(\frac{D\psi(x^*, x_0) M^2}{\zeta^2} \right) \quad \text{iterations}$$